

WHAT IS CLAIMED IS:

1 1. A method of downloading data sets by a plurality of web crawlers from among a  
2 plurality of host computers, comprising the steps of:

3 assigning a web crawler identifier to each one of the plurality of web crawlers;  
4 for each respective web crawler:

5 downloading at least one data set that includes addresses of one or more  
6 referred data sets;

7 identifying the addresses of the one or more referred data sets, wherein each  
8 identified address includes a host computer identifier;

9 for each identified address:

10 generating a representation of the host computer identifier;

11 determining a web crawler identifier to which the representation  
12 corresponds; and

13 when the determined web crawler identifier is not assigned to the  
14 respective web crawler, sending the identified address to the web crawler to which the  
15 determined web crawler identifier is assigned.

1 2. The method of claim 1, wherein

2 the plurality of web crawlers consists of  $n$  web crawlers; and

3 generating the representation includes computing a function of the host computer  
4 identifier to generate an integer value that is a member of a set of  $n$  predefined distinct values.

1 3. The method of claim 1, wherein

2 the plurality of web crawlers consists of  $n$  web crawlers; and

3 generating the representation includes computing a hash function of the host computer  
4 identifier to generate an intermediate value  $V$ , and computing  $V$  modulo  $n$ .

1 4. The method of claim 1, wherein the sending step includes:

2 determining a web crawler address for the web crawler to which the determined web  
3 crawler identifier is assigned;

4 transmitting the identified data set address to the destination web crawler at the  
5 determined web crawler address.

1 5. A method of downloading data sets by a plurality of web crawlers from among a  
2 plurality of host computers, comprising the steps of:  
3 for each respective web crawler:  
4 receiving addresses of one or more data sets from each of the plurality of web  
5 crawlers other than the respective web crawler;  
6 for each received address:  
7 determining if the address has been previously stored; and  
8 if this determination is negative, storing the address.

1 6. A web crawler system for downloading data set addresses from among a plurality of  
2 host computers, comprising:  
3 a plurality of web crawlers, wherein each web crawler has been assigned a web  
4 crawler identifier;  
5 for each respective web crawler:  
6 a main web crawler module for downloading and processing data sets stored  
7 on a plurality of host computers, the main web crawler module identifying addresses of the  
8 one or more referred data sets in the downloaded data sets, wherein each identified address  
9 includes a host computer identifier; and  
10 an address distribution module for processing the identified addresses, the  
11 address distribution module including instructions for:  
12 generating a representation of the host computer identifier, wherein the  
13 representation corresponds to one of the web crawler identifiers;  
14 determining a web crawler identifier to which the representation  
15 corresponds; and  
16 when the determined web crawler identifier is not assigned to the  
17 respective web crawler, sending the identified address to a destination web crawler  
18 comprising the web crawler to which the determined web crawler identifier is assigned.

1 7. The web crawler system of claim 6 wherein  
2 the plurality of web crawlers consists of n web crawlers; and

the address distribution module's instructions for generating the representation includes instructions for computing a hash function of the host computer identifier to generate an intermediate value V, and computing V modulo n.

8. The web crawler system of claim 6, further comprising:  
for each respective web crawler, a web crawler interface for transmitting the identified address to the destination web crawler and for receiving identified addresses from each of the plurality of web crawlers other than the respective web crawler.

9. The web crawler system of claim 6, further comprising:  
for each respective web crawler, a lookup table storing for each of the plurality of web crawler identifiers a corresponding web crawler address, said lookup table for use by the address distribution module in determining a web crawler address to which to send the identified data set address.

10. A computer program product for use in conjunction with a web crawler system wherein each web crawler is assigned a web crawler identifier, the computer program product comprising a computer readable storage medium and a computer program mechanism embedded therein, the computer program mechanism comprising:

a main web crawler module for downloading and processing data sets stored on a plurality of host computers, the main web crawler module identifying addresses of the one or more referred data sets in the downloaded data sets, wherein each identified address includes a host computer identifier; and

an address distribution module for processing the identified addresses, the address distribution module including instructions for:

generating a representation of the host computer identifier, wherein the representation corresponds to one of the web crawler identifiers;

determining a web crawler identifier to which the representation corresponds;

and

when the determined web crawler identifier is not assigned to the respective web crawler, sending the identified address to a destination web crawler comprising the web crawler to which the determined web crawler identifier is assigned.

1 11. The computer program product of claim 10, wherein:

2 the web crawler system consists of n web crawlers; and

3 the address distribution module's instructions for generating the representation  
4 includes instructions for computing a function of the host computer identifier to generate an  
5 integer value that is a member of a set of n predefined distinct values.

1 12. The computer program product of claim 10, wherein:

2 the web crawler system consists of n web crawlers; and

3 the address distribution module's instructions for generating the representation  
4 includes instructions for computing a hash function of the host computer identifier to  
5 generate an intermediate value V, and computing V modulo n.

1 13. The computer program product of claim 10, further comprising:

2 a web crawler interface for transmitting the identified address to the destination web  
3 crawler and for receiving identified addresses from each of the plurality of web crawlers other  
4 than the respective web crawler.

1 14. The computer program product of claim 10, further comprising:

2 a lookup table storing for each of the plurality of web crawler identifiers a  
3 corresponding web crawler address, said lookup table for use by the address distribution  
4 module in determining a web crawler address to which to send the identified data set address.